

## Approximate iterations for structured matrices

Wolfgang Hackbusch · Boris N. Khoromskij ·  
Eugene E. Tyrtysnikov

Received: 30 November 2005 / Revised: 16 October 2007 / Published online: 27 February 2008  
© The Author(s) 2008

**Abstract** Important matrix-valued functions  $f(A)$  are, e.g., the inverse  $A^{-1}$ , the square root  $\sqrt{A}$  and the sign function. Their evaluation for large matrices arising from pdes is not an easy task and needs techniques exploiting appropriate structures of the matrices  $A$  and  $f(A)$  (often  $f(A)$  possesses this structure only approximately). However, intermediate matrices arising during the evaluation may lose the structure of the initial matrix. This would make the computations inefficient and even infeasible. However, the main result of this paper is that an iterative fixed-point like process for the evaluation of  $f(A)$  can be transformed, under certain general assumptions, into another process which preserves the convergence rate and benefits from the underlying structure. It is shown how this result applies to matrices in a tensor format with a bounded tensor rank and to the structure of the hierarchical matrix technique. We demonstrate our results by verifying all requirements in the case of the iterative computation of  $A^{-1}$  and  $\sqrt{A}$ .

---

This work was performed during the stay of the third author at the Max-Planck-Institute for Mathematics in the Sciences (Leipzig) and also supported by the Russian Fund of Basic Research (grants 05-01-00721, 04-07-90336) and a Priority Research Grant of the Department of Mathematical Sciences of the Russian Academy of Sciences.

---

W. Hackbusch (✉) · B. N. Khoromskij  
Max-Planck-Institut für Mathematik in den Naturwissenschaften,  
Inselstr. 22-26, 04103 Leipzig, Germany  
e-mail: wh@mis.mpg.de

B. N. Khoromskij  
e-mail: bokh@mis.mpg.de

E. E. Tyrtysnikov  
Institute of Numerical Mathematics,  
Russian Academy of Sciences, Gubkina 8,  
119991 Moscow, Russia  
e-mail: tee@inm.ras.ru

# Mathematics Subject Classification (2000) 65F30 · 65F50 · 65N35 · 65F10

## 1 Introduction

We consider important matrix-valued functions  $f(A)$  as, e.g., the inverse  $A^{-1}$  and the square root  $\sqrt{A}$  [3, 10, 28–30, 41]. In particular, we are interested in evaluations of  $f(A)$  for matrices  $A$  arising from partial differential equations. Obviously, the computation of  $f(A)$  for large-scale matrices  $A$  is not an easy task. In the numerical treatment one has to avoid the full-matrix representation. Instead one should use special representations (i.e., special data structures) which, on the other hand, correspond to special properties of the argument  $A$ , of the result  $f(A)$  and of the auxiliary matrices arising during the computation process.

Examples for such a representation are Toeplitz-like structures or a sparse-matrix format. The latter format is not successful for our examples, since sparse matrices  $A$  produce results  $A^{-1}$ ,  $\sqrt{A}$ ,  $\text{sign}(A)$ , which are usually non-sparse and which, moreover, cannot be approximated by sparse matrices. This is different for the format of hierarchical matrices (cf. [19, 21–23]), the hierarchical Kronecker-tensor product (HKT) representation (cf. [24, 25, 27]) and standard Kronecker representation [36, 37].

The matrices belonging to a particular representation are characterised by a subset  $S$  of the vector space of matrices. The letter  $S$  abbreviates “structured matrices”. In the simplest case,  $A \in S$  implies  $f(A) \in S$ . If also all intermediate results belong to  $S$ , the whole computational process can be performed using the special data structures of  $S$ . The purpose of this paper is the analysis of a more complicated situation, when  $A \in S$  does not imply  $f(A) \in S$ , but  $f(A)$  has a good approximation in  $S$ . We illustrate this situation by the following example.

We consider a discrete two-dimensional Laplacian

$$A = T \otimes I + I \otimes T, \quad T = \begin{bmatrix} 2 & -1 & & \\ -1 & 2 & -1 & \\ & \ddots & \ddots & \ddots \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix}, \quad (1.1)$$

where  $T$  and the identity  $I$  are  $n \times n$  matrices. Obviously,  $A$  is a matrix of size  $n^2 \times n^2$  with a very special structure: it is *exactly* the sum of two terms, each being the Kronecker (tensor) product of two  $n \times n$  matrices. It is remarkable that  $A^{-1}$  is *approximately* of the same structure but with a greater number of terms. This number is called the tensor rank; the way the rank depends on the approximation accuracy  $\varepsilon$  and  $n$  can be seen from Table 1. Similar results for the square root of  $A$  are presented in Table 2.

We observe a logarithmic growth of the tensor rank upon  $\varepsilon$  and as well upon  $n$ . More precisely, the rank estimate  $r = O(|\log \varepsilon| \log n)$  can be proven (cf. [26]) based on approximation by exponential sums also for Kronecker products involving more than two factors (cf. [18, 24, 25]). Thus,  $A^{-1}$  and  $\sqrt{A}$  can be approximated by a matrix defined by a reasonably small number of parameters in the tensor format.

**Table 1** Tensor ranks for  $\varepsilon$ -approximations to  $A^{-1}$ 

$n$	$\varepsilon$							
	$10^{-2}$	$10^{-3}$	$10^{-4}$	$10^{-5}$	$10^{-6}$	$10^{-7}$	$10^{-8}$	$10^{-9}$
20	4	5	6	7	8	9	10	10
40	4	6	7	8	10	11	12	13
80	4	6	8	10	11	13	14	15
160	4	7	9	11	13	14	16	18
320	5	7	10	12	14	16	18	20

**Table 2** Tensor ranks for  $\varepsilon$ -approximations to  $\sqrt{A}$ 

$n$	$\varepsilon$							
	$10^{-2}$	$10^{-3}$	$10^{-4}$	$10^{-5}$	$10^{-6}$	$10^{-7}$	$10^{-8}$	$10^{-9}$
80	2	3	5	7	8	10	11	13
160	2	3	5	7	9	11	13	15
320	2	3	5	7	9	12	14	16

**Table 3** Tensor ranks for  $\varepsilon$ -approximations to  $X_k$  ( $n = 160$ )

$\varepsilon$	$k$															
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
$10^{-3}$	2	3	4	4	5	5	6	6	5	6	6	6	7	7	7	7
$10^{-6}$	2	4	7	8	8	9	10	10	11	12	12	13	14	14	13	13

So far the existence of an approximation  $\tilde{B} \approx A^{-1}$  with  $\tilde{B} \in S$  is ensured (here, the set  $S$  of structured matrices is given by sums of Kronecker products with a certain limited number of terms). It remains to design an algorithm for computing  $f(A) = A^{-1}$ . A possible choice is the Newton iteration

$$X_0 = \alpha I, \quad X_k = X_{k-1}(2I - AX_{k-1}) \quad (k = 1, 2, \dots).$$

For this iteration it can be proved that if  $0 < \alpha \leq 1/4$  then  $X_k \rightarrow A^{-1}$ , and the convergence is quadratic.

Here the important question arises, whether the *intermediate matrices*  $X_k$  belong to the subset  $S$  or can be well approximated by  $\tilde{X}_k \in S$ . In the following numerical experiment each  $X_k$  admits a suitable approximation of low tensor rank as can be seen from Table 3.

Thus, a natural idea is to substitute  $X_k$  by its approximation in the tensor format. Such a substitution is called *truncation*. Assume that the truncation is performed at every iteration. Then the following questions arise: How will this affect the convergence rate of the Newton-Schulz method? Will the convergence remain quadratic? The

answers are positive. Moreover, the same answer is valid not only for the Laplacian but typical for the truncation based on the tensor or hierarchical formats [20, 24, 25, 32].

It is worthy to mention that the truncation error may vary during iterations. It makes sense to take it not very small during the first iterations and then tighten it gradually towards the end. This strategy really helps to compress the intermediates [32].

The very idea of iterations with truncation has been already advocated in several papers, chiefly for Toeplitz-like matrices [8, 9, 31–33], rank-structured matrices [4, 5, 12, 13, 21, 27, 32] (see also [22–26]) and wavelet-based sparsification [1, 6, 12, 13]. However, the proofs provided so far only for some particular cases of structures have appeared as different individual proofs. Now we realise that many cases of previous works can be covered by a single common proof.

The main result of this paper is that an iterative fixed-point process for the evaluation of  $f(A)$  can be transformed, under certain general assumptions, into another process which preserves the convergence rate and benefits from the underlying structure. It is shown how this result applies to matrices in a tensor format with a bounded tensor rank and to the structure of the hierarchical matrix technique. We demonstrate our results by verifying all requirements in the case of the iterative computation of  $A^{-1}$  and  $\sqrt{A}$ .

In this paper we propose a general framework in which the above-mentioned results appear as particular cases. Our main results are two theorems (Sect. 2) that turn out to be both entirely general and quite elementary. Despite the latter, they do not seem to be well-known in the community of numerical analysis and structured matrices. Our results clearly amplify the role of non-linear iterative schemes in computations with structured matrices. It is especially gainful that they apply to many interesting iterative scheme and various classes of structured matrices including those already in work and those that may appear yet in application contexts.

Nevertheless, there are iterations which do not satisfy the requirements of our theorems. For instance, our theory does not apply when the success of the iteration depends on the fact that the iterates  $X_k$  stay in some sub-manifold.

The rest of the paper is organised as follows.

In Sect. 2 we consider an iteration  $X_k = \Phi_k(X_{k-1})$ , which starting with  $X_0 := \Phi_0(A)$  is assumed to converge to  $f(A)$ . The quadratic convergence is described in detail in Lemma 2.1. Next we introduce a so-called *truncation operator*  $R$  which maps into a subset  $S$  (which we call the set of “structured” elements). The combination of the iteration with the truncation operator yields the iteration with truncation  $Y_k = R(\Phi_k(Y_{k-1}))$ . In Theorem 2.2 we describe the characteristic requirements on  $R$  so that the iteration with truncation has similar convergence properties as the original iteration. The final Theorem 2.4 considers the important case that the desired result  $f(A)$  does not belong to  $S$  but is close to  $S$ .

It remains to verify that the assumptions on the truncation operator  $R$  can be satisfied in practically relevant cases. In Sect. 3 we describe a general framework which is later applied (i) to the structure used in the hierarchical matrix technique and (ii) to low Kronecker rank matrices.

Section 4 is devoted to the convergence analysis of certain matrix iterations resulting in  $A^{-1}$  and  $\sqrt{A}$  (note that the latter case is closely related to the computation of  $\text{sign}(A)$ )

[29]). In particular, the general theory from Sects. 2 and 3 ensures the quadratic convergence of the truncated Newton iterations to compute  $A^{-1}$  and  $\sqrt{A}$ .

## 2 Main result

### 2.1 Exact iteration

Let  $V$  be a normed space  $V$  and consider a function  $f : V \rightarrow V$  and  $A \in V$ . Assume that  $B := f(A)$  can be obtained by an iteration of the form

$$X_k = \Phi_k(X_{k-1}), \quad k = 1, 2, \dots, \quad (2.1)$$

where  $\Phi_k$  is a one-step operator. Further, assume that for any initial guess  $X_0$  sufficiently close to  $B$ , the process converges:

$$\lim_{k \rightarrow \infty} X_k = B. \quad (2.2)$$

If  $\Phi_k = \Phi$  does not depend on  $k$ , (2.1) represents the important *fixed-point iteration*.

**Lemma 2.1** *Let  $B$  and  $\Phi_k$  be as above and assume that there are constants  $c_\Phi, \varepsilon_\Phi > 0$  and  $\alpha > 1$  such that*

$$\begin{aligned} \|\Phi_k(X) - B\| &\leq c_\Phi \|X - B\|^\alpha \quad \text{for all } X \in V \text{ with } \|X - B\| \\ &\leq \varepsilon_\Phi \text{ and all } k \in \mathbb{N}, \end{aligned} \quad (2.3)$$

and set

$$\varepsilon := \min(\varepsilon_\Phi, 1/c), \quad c := \alpha^{-1} \sqrt[\alpha]{c_\Phi}. \quad (2.4)$$

Then (2.2) holds for any initial guess  $X_0$  satisfying  $\|X_0 - B\| < \varepsilon$ , and, moreover,

$$\|X_k - B\| \leq c^{-1} (c \|X_0 - B\|)^\alpha \quad (k = 0, 1, 2, \dots). \quad (2.5)$$

*Proof* Let  $e_k := \|X_k - B\|$ . Then, due to (2.3),

$$e_k \leq c_\Phi e_{k-1}^\alpha, \quad \text{provided that } e_{k-1} \leq \varepsilon_\Phi. \quad (2.6)$$

Because of (2.6), the inequalities  $e_{k-1} \leq \varepsilon \leq \varepsilon_\Phi$  imply  $e_k \leq c_\Phi \varepsilon^\alpha = c^{\alpha-1} \varepsilon^\alpha = \varepsilon (c\varepsilon)^{\alpha-1} \leq \varepsilon$ . Hence, all iterates stay in the  $\varepsilon$ -neighbourhood of  $B$ . (2.5) is proved by induction:

$$\begin{aligned} e_k &\stackrel{(2.6)}{\leq} c_\Phi e_{k-1}^\alpha \stackrel{\text{induction}}{=} c_\Phi \cdot \left( c^{-1} (ce_0)^{\alpha^{k-1}} \right)^\alpha \\ &\stackrel{c_\Phi = c^{\alpha-1}}{=} c^{\alpha-1} \cdot c^{-\alpha} (ce_0)^{\alpha^k} = c^{-1} (ce_0)^{\alpha^k}. \end{aligned}$$

Whenever  $e_0 < \varepsilon$ , (2.5) shows  $e_k \rightarrow 0$ . □

We remark that (2.6) together with  $e_0 \leq \varepsilon$  implies monotonicity:

$$\|X_k - B\| \leq \|X_{k-1} - B\|. \quad (2.7)$$

## 2.2 Iteration with truncation

Let  $S \subset V$  be a subset (not necessarily a subspace) considered as a class of certain structured elements (e.g., matrices of a certain data structure) and suppose that  $R : V \rightarrow S$  is an operator from  $V$  onto  $S$ . We call  $R$  a *truncation operator*. It is assumed that  $R(X) = X$  for any  $X \in S$  (i.e., all elements in  $S$  are fixed points of  $R$ ). Note that, in general,  $R$  is a non-linear mapping. The truncation of real numbers to machine numbers is a common example for  $V = \mathbb{R}$ .

Now, instead of (2.1), consider an *iterative process with truncation* defined as follows:

$$\begin{aligned} Y_0 &:= R(X_0), \\ Y_k &:= R(\Phi_k(Y_{k-1})) \quad (k = 1, 2, \dots). \end{aligned} \quad (2.8)$$

The next theorem needs the assumption that the desired result  $B := f(A)$  belongs (exactly) to the subset  $S$ . Later, in Theorem 2.4, this requirement will be relaxed.

**Theorem 2.2** *Under the premises of Lemma 2.1, assume that*

$$\|X - R(X)\| \leq c_R \|X - B\| \quad \text{for all } X \in V \text{ with } \|X - B\| \leq \varepsilon_\Phi. \quad (2.9)$$

*Then there exists  $\delta > 0$  such that the truncated iterative process (2.8) converges to  $B$  so that*

$$\|Y_k - B\| \leq c_{R\Phi} \|Y_{k-1} - B\|^\alpha \quad \text{with } c_{R\Phi} := (c_R + 1)c_\Phi \quad (k = 1, 2, \dots) \quad (2.10)$$

*for any starting value  $Y_0 = R(Y_0)$  satisfying  $\|Y_0 - B\| < \delta$ .*

*Proof* Let  $\varepsilon$  as in (2.4) and define  $Z_k := \Phi_k(Y_{k-1})$ . By (2.7) we have  $\|Z_k - B\| \leq \|Y_{k-1} - B\|$ , provided that  $\|Y_{k-1} - B\| \leq \varepsilon$ . Then

$$\|Y_k - B\| = \|(R(Z_k) - Z_k) + (Z_k - B)\| \leq (c_R + 1) \|Z_k - B\|. \quad (2.11a)$$

Assuming  $\|Y_{k-1} - B\| \leq \varepsilon$ , the inequalities  $\varepsilon \leq \varepsilon_\Phi$  and (2.3) ensure

$$\|Z_k - B\| = \|\Phi_k(Y_{k-1}) - B\| \leq c_\Phi \|Y_{k-1} - B\|^\alpha. \quad (2.11b)$$

Combining (2.11a) and (2.11b), we obtain (2.10) for any  $k$ , provided that  $\|Y_{k-1} - B\| \leq \varepsilon$ .

Similar to the proof of Lemma 2.1 and (2.7), the choice

$$\delta := \min(\varepsilon, 1/C), \quad C := \sqrt[\alpha-1]{c_{R\Phi}} \quad (2.11c)$$

guarantees that  $\|Y_0 - B\| \leq \delta$  implies  $\|Y_k - B\| \leq \delta \leq \varepsilon$  for all  $k \in \mathbb{N}$ .  $\square$

**Corollary 2.3** *Under the assumptions of Theorem 2.2, any starting value  $Y_0$  with  $\|Y_0 - B\| \leq \delta$  leads to*

$$\|Y_k - B\| \leq C^{-1} (C \|Y_0 - B\|)^{\alpha^k} \quad (k = 1, 2, \dots), \quad (2.12)$$

where  $C$  and  $\delta$  are defined in (2.11c).

### 2.3 The case of $B \notin S$

In most of the practical applications, the desired result  $B$  will not belong to the subset  $S$ , but may be close to  $S$ . The following requirement (2.14) states that  $\|B - R(B)\| \leq \varepsilon_{RB}$ . Then the iteration with truncation cannot converge to  $B$ , but it comes sufficiently close to  $B$ . In fact, in a first phase the iteration with truncation is described by (2.12) with  $C$  replaced by  $C' := \sqrt[\alpha-1]{2c_{R\Phi}}$  until it reaches the  $2\varepsilon_{RB}$ -neighbourhood of  $B$ . The quantity  $\varepsilon_{RB}$  must be sufficiently small:

$$\varepsilon_{RB} < \frac{\eta}{2}, \quad \text{where } \eta := \min\left(\varepsilon, 1/\sqrt[\alpha-1]{2c_{R\Phi}}\right) \quad (2.13)$$

with  $c_{R\Phi} = (c_R + 1)c_\Phi$  as defined above.

**Theorem 2.4** *Under the premises of Lemma 2.1, suppose*

$$\|X - R(X)\| \leq c_R \|X - B\| + \varepsilon_{RB} \quad \text{for all } X \in V \text{ with } \|X - B\| \leq \varepsilon_\Phi, \quad (2.14)$$

where  $\varepsilon_{RB}$  satisfies (2.13). Further, assume  $\|Y_0 - B\| < \eta$  and define  $Y_k$  by the iteration with truncation (2.8). Let  $m$  be the minimal  $k \in \mathbb{N}$  such that

$$\|Y_{k-1} - B\|^\alpha \leq \frac{\varepsilon_{RB}}{c_{R\Phi}}. \quad (2.15)$$

Then the errors  $\|Y_k - B\|$  strictly decrease for  $1 \leq k < m$ , while for  $k \geq m$  the iterates stagnate in a  $2\varepsilon_{RB}$ -neighbourhood of the true result:

$$\|Y_k - B\| \leq \begin{cases} 2c_{R\Phi} \|Y_{k-1} - B\|^\alpha & \text{for } k \leq m-1, \\ 2\varepsilon_{RB} & \text{for } k \geq m. \end{cases} \quad (2.16)$$

*Proof* Instead of (2.11a) we now have

$$\|Y_k - B\| \leq \|Y_k - Z_k\| + \|Z_k - B\| \leq (c_R + 1) \|Z_k - B\| + \varepsilon_{RB},$$

which obviously implies

$$\|Y_k - B\| \leq c_{R\Phi} \|Y_{k-1} - B\|^\alpha + \varepsilon_{RB}. \quad (2.17)$$

If  $k < m$ , the inequality  $\varepsilon_{RB} \leq c_{R\Phi} \|Y_{k-1} - B\|^\alpha$  holds and implies  $\|Y_k - B\| \leq 2c_{R\Phi} \|Y_{k-1} - B\|^\alpha$ . Hence, (2.10) holds with  $c_{R\Phi}$  replaced by  $2c_{R\Phi}$  giving rise to (2.12) with  $C$  replaced by  $C' := \sqrt[\alpha-1]{2c_{R\Phi}}$ . The initial error estimate  $\|Y_0 - B\| < \eta$  implies the strict decrease of  $\|Y_k - B\|$  until (2.15) holds.

If  $k = m$ , (2.17) shows  $\|Y_m - B\| \leq 2\varepsilon_{RB}$ . For  $k \geq m$ , the estimate  $c_{R\Phi} (2\varepsilon_{RB})^\alpha + \varepsilon_{RB} \leq 2\varepsilon_{RB}$  derived from (2.13) proves the second case in (2.16).  $\square$

**Corollary 2.5** *Theorems 2.2 and 2.4 can be generalised by replacing the conditions (2.9) and (2.14) with the respective inequalities*

$$\|(I - R)(X)\| \leq c_R \|X - B\|^\beta \quad (2.18)$$

and

$$\|(I - R)(X)\| \leq c_R \|X - B\|^\beta + \varepsilon_B, \quad (2.19)$$

provided that  $\alpha\beta > 1$ . Then, the order of convergence of the truncated iterative process (2.8) becomes  $\alpha\beta$ . However, all truncation operators used in this paper satisfy the conditions with  $\beta = 1$ .

Note that condition (2.9) has a clear geometrical background. If

$$R(X) := \operatorname{argmin} \{\|X - Y\| : Y \in S\}$$

is a best approximation to  $X$  in the given norm, inequality (2.9) holds with  $c_R = 1$ , since  $B \in S$ . Therefore, (2.9) with  $c_R \geq 1$  can be viewed as a *quasi-optimality condition*. If the norm is defined by a scalar product, then  $S$  is a subspace,  $R(X)$  is the orthogonal projection onto  $S$  and (2.9) is obviously fulfilled with  $c_R = 1$ .

The requirement  $\alpha > 1$  for the order of convergence implies convergence in a suitable neighbourhood of  $B$ . For linear convergence ( $\alpha = 1$ ) the additional requirement  $c_\Phi < 1$  is essential.

**Remark 2.6** In the case of  $\alpha = 1$  (i.e., linear convergence), the truncated process retains linear convergence, provided that  $(c_R + 1)c_\Phi < 1$ .

### 3 Truncation operators

Theorems 2.2 and 2.4 can be applied to various classes of structured matrices. When constructing a truncation operator for a particular class, we should take care that condition (2.9) is satisfied.



### 3.1 General framework

Next we describe a general framework which seems to cover all important cases.

**Lemma 3.1** *Let  $B = R(B)$  be fixed and assume that  $R$  is Lipschitz at  $B$ . Then the inequality (2.9) holds.*

*Proof* The Lipschitz property of  $R$  means that  $\|R(X) - R(B)\| \leq c \|X - B\|$  for some constant  $c > 0$  independent of  $X$ . The estimate

$$\|X - R(X)\|_{B=R(B)} = \|(X - B) + (R(B) - R(X))\| \leq (1 + c) \|X - B\|$$

shows (2.9) with  $c_R = 1 + c$ .

**Corollary 3.2** *Condition (2.9) is fulfilled as soon as  $B = R(B)$  and  $R$  is a bounded linear operator.*

Let  $V = \mathbb{R}^{I \times I}$  be the space of square matrices with respect to the index set  $I$  and  $S \subset V$  a subspace with a prescribed sparsity pattern  $P \subset I \times I$ , i.e.,  $X \in S$  if and only if  $X_{ij} = 0$  for all  $(i, j) \notin P$ . A familiar example of a truncation in this case is  $R(X)$  defined entry-wise by

$$R(X)_{ij} := \begin{cases} X_{ij} & \text{for } (i, j) \in P, \\ 0 & \text{for } (i, j) \notin P. \end{cases} \quad (3.1)$$

This  $R$  is linear, and hence, satisfies the hypotheses of Lemma 3.1 via Corollary 3.2.

There are only rare examples, for which  $A$  and  $B = f(A)$  can simultaneously be approximated by sparse matrices from  $S := \{X \in \mathbb{R}^{I \times I} : R(X) = X\}$ . However, it is well-known that after a discrete wavelet transform  $X \mapsto L(X) := T^{-1}XT$  one can apply a matrix compression (see [12, 13, 31, 32]). Such a matrix compression is of the form (3.1) and will be denoted by  $\Pi$  instead of  $R$ . Then, the truncation  $R$  applied to the original matrix  $X$  is the composition of the wavelet transform  $L$ , the pattern projection  $\Pi$  and the back-transformation  $L^{-1}$ :

$$R := L^{-1} \circ \Pi \circ L. \quad (3.2)$$

The same product form of  $R$  is typical as well for many other choices of  $L$  and  $\Pi$ .

In the following lemmata the operator  $\Pi$  may be non-linear.

**Lemma 3.3** *Let  $V$  and  $W$  be normed spaces and  $L : V \rightarrow W$  a bounded linear operator with a bounded inverse. Given  $B \in V$ , assume that  $\Pi : W \rightarrow W$  satisfies*

$$\|Z - \Pi(Z)\| \leq c_\Pi \|Z - L(B)\| \quad \text{for all } Z \in W \text{ with } \|L^{-1}(Z) - B\| \leq \varepsilon_\Phi. \quad (3.3)$$

*Then the truncation operator  $R$  of the form (3.2) satisfies condition (2.9) with*

$$c_R := c_\Pi \|L\| \|L^{-1}\|. \quad (3.4)$$

*Proof* Let  $Z = L(X)$ . Then, obviously,

$$\|R(X) - X\| = \left\| L^{-1}(\Pi(Z) - Z) \right\| \leq c_{\Pi} \|L^{-1}\| \|Z - L(B)\|,$$

and it remains to observe that  $\|Z - L(B)\| = \|L(X) - L(B)\| \leq \|L\| \|X - B\|$ .  $\square$

Applications of Lemma 3.3 (especially in the case of hierarchical block matrices) are facilitated by the following construction. Define a suitable system of normed spaces  $W_1, \dots, W_N$  and set

$$\begin{aligned} W &:= W_1 \times \dots \times W_N = \{H = (H_1, \dots, H_N) : H_i \in W_i\} \\ \text{with } \|H\| &= \sqrt{\sum_{i=1}^N \|H_i\|^2}. \end{aligned} \quad (3.5)$$

Let each  $W_i$  be associated with a truncation operator  $\Pi_i : W_i \rightarrow W_i$  satisfying

$$\|H_i - \Pi_i(H_i)\| \leq c_i \|H_i - Z_i\| \quad \text{for all } H_i \in W_i \text{ and } 1 \leq i \leq N, \quad (3.6)$$

where  $Z_i \in W_i$  are some fixed elements.

**Lemma 3.4** *Let  $W$  be the normed space from (3.5) and let the truncation operators  $\Pi_i$  satisfy (3.6), where the elements  $Z_i \in W_i$  are defined by*

$$L(B) = (Z_1, \dots, Z_N).$$

*The product of the truncation operators  $\Pi_i$  defines  $\Pi : W \rightarrow W$  via*

$$\Pi(H) := (\Pi_1(H_1), \dots, \Pi_N(H_N)) \quad \text{for } H = (H_1, \dots, H_N), \quad H_i \in W_i.$$

*Then  $R$  from (3.2) satisfies (2.9).*

*Proof* Let  $L(X) = H = (H_1, \dots, H_N)$ . Then, according to the definitions of  $L$  and  $\Pi$ ,

$$\|H - \Pi(H)\| \leq \sqrt{\sum_{i=1}^N c_i^2 \|H_i - Z_i\|^2} \leq \left( \max_{1 \leq i \leq N} c_i \right) \sqrt{\sum_{i=1}^N \|H_i - Z_i\|^2},$$

which proves (3.3) and allows us to use Lemma 3.3.  $\square$

An important example of  $\Pi$  in the case of a matrix space  $W$  is given by optimal low-rank approximations.

**Lemma 3.5** *Let  $W$  be a normed space of all matrices of a fixed size and let  $S \subset W$  consist of all matrices whose rank does not exceed  $r$ . Then for any  $H \in W$  there exists a matrix  $T \in S$  such that  $\|H - T\| = \min_{\text{rank } Z \leq r} \|H - Z\|$ .*

*Proof* Consider a minimising sequence  $Z_k \in S$ , i.e.,  $\lim_{k \rightarrow \infty} \|H - Z_k\| = \delta := \inf_{\text{rank } Z \leq r} \|H - Z\|$ . Obviously, the sequence  $Z_k$  is bounded. Therefore, a convergent subsequence  $Z_{k_i} \rightarrow T$  exists. Its limit satisfies  $\|H - T\| = \delta$ .

The assertion  $T \in S$  is due to the fact that a matrix of rank equal to  $p > r$  possesses a vicinity wherein any matrix is of rank  $\geq p$ .  $\square$

The optimal approximant  $T$  is not necessarily unique. For the mathematical definition of  $\Pi(H)$  we choose any of the optimal approximants. In practice, the result depends on the implementation.

**Corollary 3.6** *For any norm, the optimal truncation operator  $\Pi$  defined in Lemma 3.3 satisfies (3.3) with  $c_\Pi = 1$ .*

*Proof* In the given norm, no matrix in  $S$  can be closer to  $H$  than  $\Pi(H)$ .  $\square$

Matrix theory provides well-developed tools for the construction of low-rank approximations in the case of any unitarily invariant norm. For an arbitrary matrix  $H \in W$ , denote its singular values by  $\sigma_1(H) \geq \sigma_2(H) \geq \dots$  and let  $\Sigma(H) := \text{diag}\{\sigma_1(H), \sigma_2(H), \dots\}$ . Let  $\Sigma_r(H)$  be obtained from  $\Sigma(H)$  by retaining all  $\sigma_k(H)$  for  $1 \leq k \leq r$  and changing the other entries into zeroes. Let  $H = Q_1 \Sigma(H) Q_2$  be the singular value decomposition of  $H$  (with unitary  $Q_1$  and  $Q_2$ ). Then

$$\Pi(H) := Q_1 \Sigma_r(H) Q_2 \quad (3.7)$$

is the best possible approximant to  $H$  in the set  $S$  of matrices of rank  $\leq r$ , where the norm is arbitrary but unitarily invariant. It can be readily deduced from the Mirsky theorem (cf. [7, 35]) claiming that

$$\|\Sigma(H) - \Sigma(Z)\| \leq \|H - Z\| \quad (3.8)$$

for all matrices  $H$  and  $Z$  of the same size and any unitarily invariant norm. If  $Z \in S$ , then, clearly,  $\sigma_i(Z) = 0$  for  $i \geq r + 1$ . Using this together with the monotonicity of unitarily invariant norms (cf. [35]), we obtain

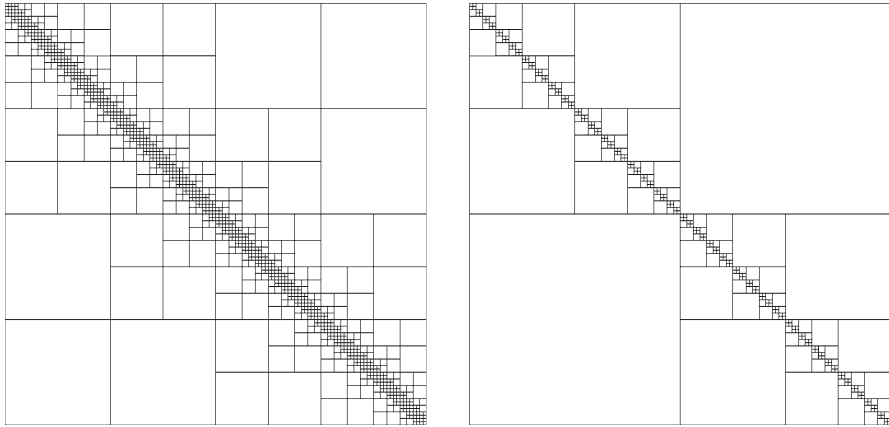
$$\|H - \Pi(H)\| = \|\Sigma(H) - \Sigma_r(H)\| \leq \|\Sigma(H) - \Sigma(Z)\|,$$

and, due to the Mirsky theorem, the latter norm is estimated from above by  $\|H - Z\|$ .

For the most familiar unitarily invariant norms such as the spectral and the Frobenius norm, the above facts can be established through simpler arguments. In particular, it is well-known that

$$\min_{\text{rank } Z \leq r} \|H - Z\|_2 = \sigma_{r+1}(H), \quad \min_{\text{rank } Z \leq r} \|H - Z\|_F = \sqrt{\sum_{i \geq r+1} \sigma_i^2(H)}.$$

Thus, the truncation property (2.9) is easy to achieve when a best approximation element is existing. Sometimes (e.g., for three-way approximations of bounded tensor



**Fig. 1** Standard and weakly admissible  $\mathcal{H}$ -partitionings

rank) this is not the case. Nevertheless, all cases are supported by Theorem 2.4 as we can always capitalise on a quasi-optimal construction as follows.

Let  $\delta(H) = \inf_{T \in \mathcal{S}} \|H - T\|$ . For a given fixed  $\varepsilon > 0$ , let  $\Pi(H)$  denote an  $\varepsilon$ -optimal approximation to  $H$  in the sense that

$$\delta(H) \leq \|H - \Pi(H)\| \leq \delta(H) + \varepsilon.$$

**Lemma 3.7** *If  $\Pi(H)$  is defined as an  $\varepsilon$ -optimal approximation to  $H$  on  $\mathcal{S}$ , then*

$$\|H - \Pi(H)\| \leq \|H - Z\| + \varepsilon \quad \text{for any } Z \in \mathcal{S}. \quad (3.9)$$

*Proof* Use  $\|H - \Pi(H)\| - \|H - Z\| \leq (\delta(H) + \varepsilon) - \delta(H) = \varepsilon$ .  $\square$

In the next sections, we discuss some details of the construction of  $L$  and  $\Pi$  for hierarchical block matrices and matrices in the tensor format.

Other useful applications of the same general framework are Toeplitz-like matrices, where  $L(X) := PX - XQ$  for some specially chosen fixed matrices  $P$  and  $Q$  (cf. [9, 31, 33]).

### 3.2 Application to hierarchical block matrices

Let  $V$  be the space  $\mathbb{R}^{n \times n}$  of  $n \times n$  matrices. Consider a block decomposition as depicted in Fig. 1. Let  $N$  be the number of matrix blocks. Then each matrix block belongs to a certain matrix space  $W_i$  ( $1 \leq i \leq N$ ). Given  $X \in V$ , let  $L_i(X) \in W_i$  be the  $i$ th block. The space  $W$  is defined according to (3.5).

The above-considered operator  $L : V \rightarrow W$  maps a matrix  $X$  into the  $N$ -tuple of matrix-blocks:

$$L(X) := (L_1(X), \dots, L_N(X)).$$

If the Frobenius norm is used on the spaces  $V$  and  $W_1, \dots, W_N$ , the norm induced on  $W$  is again the Frobenius norm. Obviously,  $\|X\|_F = \|L(X)\|_F$  holds. Hence, the inverse  $L^{-1}$  exists and satisfies

$$\|L\| = \|L^{-1}\| = 1.$$

Fix a positive integer  $r$  and let  $S_i \subset W_i$  be the subset of matrices of rank  $\leq r$ . Define  $S$  as the Cartesian product

$$S = S_1 \times \dots \times S_N \subset W$$

and let  $\Pi_i : W_i \rightarrow S_i$  be of the form (3.7) involving the singular value decomposition of the matrix block  $W_i$ . Defining  $\Pi : W \rightarrow S$  as in Lemma 3.4 and using Lemma 3.5, we can apply Theorem 2.2 to  $R = L^{-1} \circ \Pi \circ L$ .

Note that exactly this kind of truncation is used in the theory of hierarchical block matrices (cf. [21, 22, 38, 39]) and even in some early implementations (cf. [16]).

Initially, the main purpose of the rank truncation was the reduction of storage and of the matrix-by-vector complexity. In the sequel, it was shown that with an appropriate block decomposition the hierarchical matrix structure supports all matrix operations and therefore allows to compute various matrix functions  $f(A)$  of  $A \in S \subset V$ , where  $B := f(A)$  is known to be close to  $S$  (e.g., for  $f(A) = A^{-1}$  compare [2, 15], and for  $f(A) = \text{sign}(A)$  see [14, 25]). In spite of the observation that these computations are efficient and robust, the rigorous analysis of the intermediate truncation errors was incomplete. Our results now suggest some general framework for such an analysis of basic iterative algorithms.

Finally we remark that sometimes the optimal truncation is replaced by an approximate or heuristic one which is cheaper to compute (e.g., by cross approximation techniques, see [17, 40]). However, the rigorous analysis of such kind of quasi-optimal truncation procedures is beyond the scope of our paper.

### 3.3 Application to tensor approximations

Let  $V_1 = \mathbb{R}^{p \times q}$  and  $V_2 = \mathbb{R}^{r \times s}$ , while  $V = \mathbb{R}^{pr \times qs}$  for some integers  $p, q, r, s$ . The Kronecker product is a mapping from  $V_1 \times V_2$  into  $V$ . For  $A \in V_1$  and  $B \in V_2$ , the

Kronecker product  $A \times B$  is defined by the block matrix 
$$\begin{bmatrix} a_{11}B & a_{21}B & \dots \\ a_{12}B & a_{22}B & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \in V.$$
 We

say that a matrix  $M \in V$  has a Kronecker rank  $\leq k$ , if there is a representation

$$M = \sum_{v=1}^{\ell} A_v \times B_v \quad \text{with } A_v \in V_1, B_v \in V_2, \text{ and } \ell \leq k. \quad (3.10)$$

We define the subset of structured matrices  $S$  by the set of all matrices of Kronecker rank  $\leq k$ . If  $k$  is not too large, this is an interesting representation since matrices of the large size  $pr \times qs$  can be described by matrices  $A_v, B_v$  of relatively small size.

As described, e.g., in [27], there is a simple isomorphism  $L$  from  $V = \mathbb{R}^{pr \times qs}$  to  $\mathbb{R}^{pq \times rs}$  such that the representation (3.10) of  $M \in S \subset V = \mathbb{R}^{pr \times qs}$  is equivalent to  $\text{rank}(L(M)) \leq k$ . Therefore, we obtain the situation of Lemma 3.5 with  $W := \Psi(V) = \mathbb{R}^{pq \times rs}$ . The truncation operator is again of the form  $R = L^{-1} \circ \Pi \circ L$ , where  $\Pi$  is the optimal SVD-based truncation or an appropriate substitute.

The framework of this paper can be applied also to the (multi-linear) tensor representation (3.10) where the number of factors is greater than 2. In this case the truncation procedures are not so well developed; however, some algorithms are available and claimed to be efficient in particular applications (mostly for data analysis in chemometrics, psychometrics, etc.; cf. [11]).

## 4 Examples of approximate iterations

We will consider iterative schemes to compute the matrix-valued functions  $f(A) = A^{-1}$  and  $f(A) = \sqrt{A}$ . The common feature of the considered iterative schemes is that they have locally quadratic convergence and require only matrix-matrix products in each step of the iteration. We prove that our general results can be applied in the case of hierarchical matrices, Kronecker products or mixed hierarchical Kronecker-product formats to compute  $A^{-1}$  (cf. Sect. 4.1) and  $\sqrt{A}$  (cf. Sect. 4.2.2).

On the other hand, our convergence theory for iteration with truncation does not apply, in general, to the case of Newton-type iterative schemes in a subspace.

### 4.1 Newton iteration for calculating $A^{-1}$

Let  $V = \mathbb{C}^{n \times n}$  and  $A \in V$  a regular matrix. The Newton method applied to the equation  $\Psi(X) := A - X^{-1} = 0$  yields the iteration

$$X_k := X_{k-1}(2I - AX_{k-1}) \quad (k = 1, 2, \dots), \quad (4.1)$$

which is also named Schulz iteration (cf. [34]). This corresponds to the formulation (2.1) with

$$\Phi_k(X) := \Phi(X) := X(2I - AX).$$

The Newton method is known to have locally quadratic order of convergence (i.e.,  $\alpha = 2$  in (2.3)). Let  $E_k := I - AX_k$  denote the error. Using  $X_k = X_{k-1}(I + E_{k-1})$  we obtain

$$E_k = I - AX_{k-1}(I + E_{k-1}) = I - (I - E_{k-1})(I + E_{k-1}) = E_{k-1}^2. \quad (4.2)$$

Applying (4.2) recursively, we find that

$$E_k = E_0^{2^k} \quad (k = 1, 2, \dots) \quad (4.3)$$

and conclude

$$A^{-1} - X_k = A^{-1} E_k = A^{-1} E_0^{2^k} = X_0(I - E_0)^{-1} E_0^{2^k}.$$

Hence, the iteration converges quadratically for all starting values  $X_0$  with  $\rho(E_0) < 1$ , where  $\rho$  is the spectral radius. Finally, Eq. (4.2) implies

$$A^{-1} - X_k = A^{-1} E_k = (A^{-1} - X_{k-1})A(A^{-1} - X_{k-1}),$$

which proves (2.3) with  $\alpha = 2$  and  $c_\Phi = \|A\|$ .

Now Theorem 2.4 can be applied with a proper choice of the subset  $S$  and of the truncation operator  $R$ .

## 4.2 Newton iteration for the calculation of $\sqrt{A}$

### 4.2.1 Non-constrained Newton iteration

We apply the Newton method to the equation  $\Psi(X) := A - X^2 = 0$ . Abbreviating the correction by  $\Delta_k := X_k - X_{k-1}$ , we obtain the iteration

$$X_0 \in V, \quad X_{k-1}\Delta_k + \Delta_k X_{k-1} = A - X_{k-1}^2 \quad (k = 1, 2, \dots), \quad (4.4)$$

corresponding to the choice  $\Phi_k(X) := \Phi(X)$ , where  $\Phi(X)$  solves the matrix equation

$$X(\Phi(X) - X) + (\Phi(X) - X)X = A - X^2. \quad (4.5)$$

A simple calculation shows that the latter equation implies (with the substitution  $A = B^2$ )

$$X(\Phi(X) - B) + XB - X^2 + (\Phi(X) - B)X + BX - X^2 = B^2 - X^2,$$

which leads to the matrix Lyapunov equation with respect to  $Y = \Phi(X) - B$ ,

$$XY + YX = (B - X)^2.$$

Making use of the solution operator for the Lyapunov equation [14] (and assuming that  $X = X^\top$  is positive definite), we arrive at the norm estimate

$$\|\Phi(X) - B\| = \left\| \int_0^\infty e^{-tX} (B - X)^2 e^{-tX} dt \right\| \leq C \|B - X\|^2.$$

This proves relation (2.3) with  $\alpha = 2$ . Hence, Theorem 2.4 applies to the truncated version of the non-linear iteration (4.4).

#### 4.2.2 Newton iteration in the subspace

Let  $A$  be diagonalisable, i.e.,  $A = T^{-1}D_AT$  for some  $T \in V$  and a non-negative diagonal matrix  $D_A$ . This gives rise to the subspace

$$V_T := \{M \in \mathbb{R}^{n \times n} : M = T^{-1}DT, D \text{ is diagonal}\} \subset V. \quad (4.6)$$

Note that  $A \in V_T$  and that all matrices from  $V_T$  commute.

We reconsider iteration (4.4) under the assumption  $X_0 \in V_T$  (this is trivially satisfied for all multiples  $X_0 = a_0A$ ). Next, it is easy to see that all iterates  $X_k$  of (4.4) belong to  $V_T$ . In particular,  $X \in V_T$  implies  $\Phi(X) \in V_T$  and the left-hand side in (4.5) can be simplified to  $2X\Phi(X) - 2X^2$ . Hence we obtain the iteration

$$X_0 = a_0A, \quad X_k := \frac{1}{2} \left( X_{k-1} + X_{k-1}^{-1}A \right) \quad (k = 1, 2, \dots), \quad (4.7)$$

where  $a_0 > 0$  is the given constant. This corresponds to the formulation (2.1) with

$$\Phi_k(X) := \Phi(X) := \frac{1}{2}(X + X^{-1}A).$$

Note that newly defined  $\Phi$  is different from  $\Phi$  in (4.5), but both coincide on  $V_T$ . In particular for starting values  $X_0 \in V_T$  both exact iterations yield the same  $X_k$ . Hence, the convergence analysis of Sect. 4.2.1 implies the same kind of convergence for the iteration (4.7).

Unfortunately, the iteration (4.7) is proved to be numerically stable only under quite restrictive assumptions [28]. As is shown in [29], useful stable versions for the iteration (4.7) are related to the computation of the sign function (cf. [30]):

$$Y_0 = A, \quad Z_0 = I, \\ Y_{k+1} = \frac{1}{2} \left( Y_k + Z_k^{-1} \right), \quad Z_{k+1} = \frac{1}{2} \left( Z_k + Y_k^{-1} \right), \quad k = 0, 1, \dots; \quad (4.8)$$

$$Y_0 = A, \quad Z_0 = I, \\ Y_{k+1} = \frac{1}{2} Y_k (3I - Z_k Y_k), \quad Z_{k+1} = \frac{1}{2} (3I - Z_k Y_k) Z_k, \quad k = 0, 1, \dots \quad (4.9)$$

For both cases it is known that, under appropriate normalisation of  $A$ ,

$$Y_k \rightarrow A^{1/2}, \quad Z_k \rightarrow A^{-1/2}$$

with a quadratic convergence, and (4.9) takes an obvious advantage of involving only matrix–matrix multiplications. We should note, all the same, that a rigorous analysis of truncation for the iterations (4.8) and (4.9) remains an open problem, for the truncation may cause the iterates to leave the subspace of matrices commuting with  $A$ .



**Table 4** Approximate iterations (4.9) for  $\sqrt{A}$ 

Residual error	$10^{-3}$	$10^{-4}$	$10^{-5}$	$10^{-6}$	$10^{-7}$	$10^{-8}$	$10^{-9}$	$10^{-10}$
Number of iterations	9	11	12	12	13	13	13	14
Tensor rank	3	4	6	7	8	9	9	10

In the general case, we are not aware of any convenient means to keep approximate iterations inside the subspace. Despite that, numerical results still demonstrate the usefulness of truncation.

Table 4 shows the behaviour of iterations (4.9) with truncation onto the set of matrices with a bounded tensor rank. The results are obtained for the discrete 2D Laplacian (1.1) with the 1D grid size  $n = 20$ . Thus,  $A$  is a matrix of order  $n^2 = 400$  with tensor rank equal to 2. Note that the initial matrix is first normalised to have the unit Frobenius norm. We start truncation with  $\varepsilon_1 = \varepsilon_0$ , where  $\varepsilon_0$  is the prescribed residue error, and then tighten it by half at every iteration step  $k$ :

$$\varepsilon_k = \varepsilon_0 / 2^{k-1}.$$

We quit approximate iterations as soon as the residual error enjoys the stopping criterion

$$\frac{\|A - X_k^2\|_F}{\|A\|_F} \leq \varepsilon_0 \quad (4.10)$$

and then perform the final truncation with the given  $\varepsilon_0$  and recalculate the residual error. In the end we report on the iteration number for the stopping criterion (4.10) to hold and the tensor rank of the approximations to  $A^{1/2}$  and  $A^{-1/2}$ .

## 5 Concluding remarks

In this paper we proposed a unified framework for the analysis of *iterations with truncation*. The advantage of these approximate iterations is that they preserve the data-sparse structure of the intermediate matrices. The main result is that an iterative process for the evaluation of  $f(A)$  can be transformed, under very general assumptions, into an implementable process which preserves the convergence rate and benefits from the underlying structure during the iterations. It is shown how this result applies to matrices in the tensor format with a bounded tensor rank and to the hierarchical matrices (with a bounded rank of the blocks).

**Acknowledgements** The authors are grateful to all the referees for the helpful and stimulating comments.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution Noncommercial License which permits any noncommercial use, distribution, and reproduction in any medium, provided the original author(s) and source are credited.

## References

1. Alpert, B., Beylkin, G., Gines, D., Vozovoi, L.: Adaptive solution of partial differential equations in multiwavelet bases. *J. Comput. Phys.* **182**, 149–190 (2002)
2. Bebendorf, M., Hackbusch, W.: Existence of  $\mathcal{H}$ -matrix approximants to the inverse FE-matrix of elliptic operators with  $L^\infty$ -coefficients. *Numer. Math.* **95**, 1–28 (2003)
3. Beylkin, G., Coult, N., Mohlenkamp, M.J.: Fast spectral projection algorithms for density-matrix computations. *J. Comput. Phys.* **152**, 32–54 (1999)
4. Beylkin, G., Mohlenkamp, M.M.: Numerical operator calculus in higher dimensions. *Proc. Natl. Acad. Sci. USA* **99**, 10246–10251 (2002)
5. Beylkin, G., Mohlenkamp, M.M.: Algorithms for numerical analysis in high dimensions. *SIAM J. Sci. Comput.* **26**, 2133–2159 (2005)
6. Beylkin, G., Sandberg, K.: Wave propagation using bases for bandlimited functions. *Wave Motion* **41**, 263–291 (2005)
7. Bhatia, R.: *Matrix Analysis*. Springer, New York (1996)
8. Bini, D.A., Tyrtyshnikov, E.E., Yalamov, P. (eds.): *Structured Matrices: Recent Developments in Theory and Computation*. Advances in Computation. Nova Science, Huntington (2001)
9. Bini, D.A., Meini, B.: Solving block banded block Toeplitz systems with structured blocks: algorithms and applications. In: Bini, D.A., Tyrtyshnikov, E.E., Yalamov, P. (eds.) *Structured Matrices: Recent Developments in Theory and Computation*. Advances in Computation. Nova Science, Huntington (2001)
10. Byers, R., He, C., Mehrmann, V.: The matrix sign function method and the computation of invariant subspaces. *SIMAX* **18**, 615–632 (1997)
11. De Lathauwer, L., De Moor, B., Vandewalle, J.: A multilinear singular value decomposition. *SIAM J. Matrix Anal. Appl.* **21**, 1253–1278 (2000)
12. Ford, J.M., Tyrtyshnikov, E.E.: Combining Kronecker product approximation with discrete wavelet transforms to solve dense, function-related systems. *SIAM J. Sci. Comp.* **25**, 961–981 (2003)
13. Ford, J.M., Oseledets, I.V., Tyrtyshnikov, E.E.: Matrix approximations and solvers using tensor products and non-standard wavelet transforms related to irregular grids. *Russ. J. Numer. Math. Math. Model.* **19**(2), 185–204 (2004)
14. Gavrilyuk, I.P., Hackbusch, W., Khoromskij, B.N.: Data-sparse approximation to a class of operator-valued functions. *Math. Comp.* **74**, 681–708 (2005)
15. Gavrilyuk, I.P., Hackbusch, W., Khoromskij, B.N.: Tensor-product approximation to elliptic and parabolic solution operators in higher dimensions. *Computing* **74**, 131–157 (2005)
16. Goreinov, S.A., Tyrtyshnikov, E.E., Yereimin, A.Y.: Matrix-free iteration solution strategies for large dense linear systems. *Numer. Linear Algebra Appl.* **4**, 1–22 (1996)
17. Goreinov, S.A., Tyrtyshnikov, E.E., Zamarashkin, N.L.: A theory of pseudo-skeleton approximations. *Linear Algebra Appl.* **261**, 1–21 (1997)
18. Grasedyck, L.: Existence and computation of a low Kronecker-rank approximation to the solution of a tensor system with tensor right-hand side. *Computing* **72**, 247–265 (2004)
19. Grasedyck, L., Hackbusch, W.: Construction and arithmetics of  $\mathcal{H}$ -matrices. *Computing* **70**, 295–334 (2003)
20. Grasedyck, L., Hackbusch, W., Khoromskij, B.N.: Solution of large scale algebraic matrix Riccati equations by use of hierarchical matrices. *Computing* **70**, 121–165 (2003)
21. Hackbusch, W.: A sparse matrix arithmetic based on  $\mathcal{H}$ -matrices. I. Introduction to  $\mathcal{H}$ -matrices. *Computing* **62**, 89–108 (1999)
22. Hackbusch, W., Khoromskij, B.N.: A sparse  $\mathcal{H}$ -matrix arithmetic. II. Application to multi-dimensional problems. *Computing* **64**, 21–47 (2000)
23. Hackbusch, W., Khoromskij, B.N.: A sparse  $\mathcal{H}$ -matrix arithmetic: General complexity estimates. *J. Comp. Appl. Math.* **125**, 479–501 (2000)
24. Hackbusch, W., Khoromskij, B.N.: Low-rank Kronecker product approximation to multi-dimensional nonlocal operators. Part I. Separable approximation of multi-variate functions. *Computing* **76**, 177–202 (2006)
25. Hackbusch, W., Khoromskij, B.N.: Low-rank Kronecker product approximation to multi-dimensional nonlocal operators. Part II. HKT representations of certain operators. *Computing* **76**, 203–225 (2006)
26. Hackbusch, W., Khoromskij, B.N., Kriemann, R.: Hierarchical matrices based on a weak admissibility criterion. *Computing* **73**, 207–243 (2004)

27. Hackbusch, W., Khoromskij, B.N., Tyrtshnikov, E.E.: Hierarchical Kronecker tensor-product approximations. *J. Numer. Math.* **13**, 119–156 (2005)
28. Higham, N.J.: Newton's method for the matrix square root. *Math. Comput.* **46**, 537–549 (1986)
29. Higham, N.J.: Stable iterations for the matrix square root. *Numer. Algorithms* **15**, 227–242 (1997)
30. Kenney, C.S., Laub, A.J.: The matrix sign function. *IEEE Trans. Automat. Control* **40**, 1330–1348 (1995)
31. Olshevsky, V., Oseledets, I., Tyrtshnikov, E.E.: Tensor properties of multilevel Toeplitz and related matrices. *Linear Algebra Appl.* **412**, 1–21 (2006)
32. Oseledets, I.V., Tyrtshnikov, E.E.: Approximate inversion of matrices in the process of solving a hypersingular integral equation. *Comp. Math. Math. Phys.* **45**(2), 302–313 (2005) (translated from *JVM i MF* **45**(2), 315–326 (2005))
33. Pan, V.Y., Rami, Y.: Newton's iteration for the inversion of structured matrices. In: Bini, D.A., Tyrtshnikov, E.E., Yalamov, P. (eds.) *Structured Matrices: Recent Developments in Theory and Computation*. Advances in Computation, pp. 79–90. Nova Science, Huntington (2001)
34. Schulz, G.: Iterative Berechnung der reziproken Matrix. *ZAMM* **13**, 57–59 (1933)
35. Stewart, G.W., Sun, J.: *Matrix Perturbation Theory*. Academic Press, San Diego (1990)
36. Tyrtshnikov, E.E.: Tensor approximations of matrices generated by asymptotically smooth functions. *Sbornik: Mathematics* **194**(5–6), 941–954 (2003) (translated from *Mat. Sb.* **194**(6), 146–160 (2003))
37. Tyrtshnikov, E.E.: Kronecker-product approximations for some function-related matrices. *Linear Algebra Appl.* **379**, 423–437 (2004)
38. Tyrtshnikov, E.E.: Mosaic ranks and skeletons. In: *Numerical Analysis and Its Applications. Proceedings of WNAA-96. Lecture Notes in Computer Science*, vol. 1196, pp. 505–516. Springer, Berlin (1996)
39. Tyrtshnikov, E.E.: Mosaic-skeleton approximations. *Calcolo* **33**, 47–57 (1996)
40. Tyrtshnikov, E.E.: Incomplete cross approximation in the mosaic-skeleton method. *Computing* **64**, 367–380 (2000)
41. Wilkinson, J.H.: *Rounding Errors in Algebraic Processes*. Prentice Hall, Englewood Cliffs, NJ (1963)